

Neural Computing Research Group
Dept of Computer Science & Applied Mathematics
Aston University
Birmingham B4 7ET
United Kingdom
Tel: +44 (0)121 333 4631
Fax: +44 (0)121 333 4586
<http://www.ncrg.aston.ac.uk/>

Sparse Online Gaussian Processes

Lehel Csató

csatol@aston.ac.uk

Manfred Opper

opperm@aston.ac.uk

Technical Report NCRG/2001/014

October 9, 2002

*Accepted in Neural Computation
Minor corrections included^a*

^aThe authors acknowledge reader feedbacks

Abstract

We develop an approach for sparse representations of Gaussian Process (GP) models (which are Bayesian types of kernel machines) in order to overcome their limitations for large data sets. The method is based on a combination of a Bayesian online algorithm together with a sequential construction of a relevant subsample of the data which fully specifies the prediction of the GP model. By using an appealing parametrisation and projection techniques that use the RKHS norm, recursions for the effective parameters and a sparse Gaussian approximation of the posterior process are obtained. This allows both for a propagation of predictions as well as of Bayesian error measures. The significance and robustness of our approach is demonstrated on a variety of experiments.

1 Introduction

Gaussian processes (GP) (Bernardo and Smith 1994; Williams and Rasmussen 1996) provide promising Bayesian tools for modeling real-world statistical problems. Like other *kernel* based methods, such as Support Vector Machines (SVMs) (Vapnik 1995), they combine a high flexibility of the model by working in often ∞ dimensional feature spaces with the simplicity that all operations are “kernelised” – performed in the lower dimensional input space utilising positive definite kernels.

An important advantage of GPs over other non-Bayesian models is the explicit probabilistic formulation of the model. This allows the modeller to assess the uncertainty of the predictions by providing *Bayesian* confidence intervals (for regression) or posterior class probabilities (for classification). It also opens the possibility to treat a variety of other nonstandard data models (e.g. quantum inverse statistics (Lemm et al. 2000), wind-fields (Evans et al. 2000; Berliner et al. 2000)) using a kernel method.

GPs are non-parametric in the sense that the “parameters” to be learnt are *functions* f_x of a usually continuous input variable $x \in \mathbb{R}^d$. The value f_x is used as a latent variable in a likelihood $P(y|f_x, x)$ which denotes the probability of an observable output variable y given the input x . The *a priori* assumption on the statistics of f is that of a Gaussian process: any finite collection of random variables f_i is jointly Gaussian. Hence, one must specify the prior means and the prior covariance function of the variables f_x . The latter is called the *kernel* $K_0(x, x') = \text{Cov}(y, y')$ (Vapnik 1995; Kimeldorf and Wahba 1971). Thus, if a zero-mean GP is assumed, the kernel K_0 fully specifies the entire prior information about the model. Based on a set of input-output observations (x_n, y_n) ($n = 1, \dots, N$) the Bayesian approach computes the posterior distribution of the process f_x using the prior and the likelihood (Williams 1999; Williams and Rasmussen 1996; Gibbs and MacKay 1999).

A straightforward application of this simple appealing idea is impeded by two major obstacles: non-Gaussianity of the posterior process and the size of the kernel matrix $K_0(x_i, x_j)$. A first obvious problem stems from the fact that the posterior process is usually non-Gaussian (except when the likelihood itself is Gaussian in the f_x). Hence, in many important cases its analytical form precludes an exact evaluation of the multidimensional integrals that occur in posterior averages. Nevertheless, various methods have been introduced to approximate these averages. A variety of such methods may be understood as approximations of the non-Gaussian posterior process by a Gaussian one (Jaakkola and Haussler 1999; Seeger 2000), for instance in (Williams and Barber 1998) the posterior mean is replaced by the posterior maximum (MAP) and information about the fluctuations are derived by a quadratic expansion around this maximum. The computation of these approximations, which become exact for regression with Gaussian noise, require the solution of a system of coupled nonlinear equations of the size equal to the number of data-points. The second obstacle which prevents GPs from being applied to large datasets is that the matrix which couples these equations is typically not sparse.

Hence, the development of good sparse approximations are of major importance. Such approximations aim at performing the most time consuming matrix operations (inversions or diagonalisations) only on a representative *subset* of the training data. In this way, the computational time is reduced from $\mathcal{O}(N^3)$ to $\mathcal{O}(Np^2)$; where N is the total size of the training data and p is the size of the representative set. The memory requirement is $\mathcal{O}(p^2)$ as opposed to $\mathcal{O}(N^2)$. So far, a variety of sparsity techniques (Smola and Schölkopf 2000; Williams and Seeger 2001) for *batch* training of GPs have been proposed. This paper presents a new approach which combines the idea of a sparse representation with an on-line algorithm that allows for a speedup of the GP training by sweeping through a dataset only once. A different sparse approximation which also allows for an on-line processing was recently introduced by (Tresp 2000). It is based on combining predictions of models trained on smaller data subsets and needs an additional query set of inputs.

Central to our approach are exact expressions for the posterior means $\langle f_x \rangle_t$ and the posterior covariance $K_t(x, x')$ (subscripts denote the number of data points) which are derived in section 2. Although both quantities are continuous functions, they can be represented as finite linear (or respective bilinear) combinations of kernels $K_0(x, x_i)$ evaluated at the training inputs x_i (Csató et al. 2000). Using sequential projections of the posterior process on the manifold of Gaussian processes, we obtain approximate recursions for the effective parameters of these representations. Since the size of representations grows with the number of training data, we use a second type of projection to extract a smaller subset of input data (reminiscent of the “support vectors” (Vapnik 1995) or “relevance vectors” of (Tipping 2000)). This subset builds up a sparse representation of the posterior process on which all predictions of the trained GP model rely. Our approach is related to the one introduced in Wahba (1990) ch. 7. While we use the same measure for projection, we do not fix the set of basis vectors from the beginning, but decide on-line which inputs to keep.

2 Online Learning with Gaussian Processes

In Bayesian learning, all information about the parameters that we wish to infer is encoded in probability distributions (Bernardo and Smith 1994). In the GP framework, the parameters are functions and the GP priors specify a Gaussian distribution over a function space. The posterior process is entirely specified by all its finite dimensional marginals. Hence, let $\mathbf{f} = \{f(\mathbf{x}_1), \dots, f(\mathbf{x}_M)\}$ be a set of function values such that $\mathbf{f}_{\mathcal{D}} \subseteq \mathbf{f}$, where $\mathbf{f}_{\mathcal{D}}$ is the set of $f(\mathbf{x}_i) = f_i$ with \mathbf{x}_i in the observed set of inputs, we compute the posterior distribution using the data likelihood together with the prior $p_0(\mathbf{f})$ as

$$p_{\text{post}}(\mathbf{f}) = \frac{P(\mathcal{D}|\mathbf{f})p_0(\mathbf{f})}{\langle P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) \rangle_0}, \quad (1)$$

where $\langle P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) \rangle_0$ is the average of the likelihood with respect to the prior GP (GP at time 0). This form of the posterior distribution can be used to express posterior expectations as typically high dimensional integrals. For prediction, one is especially interested in expectations of functions of the process at inputs, which are not contained in the training set. At first glance, one might assume that every prediction on a novel input would require the computation of a new integral. Even if we had good methods for approximate integration, this would make predictions at new inputs a rather tedious task. Luckily, the following lemma shows that simple but important predictive quantities like the posterior mean and the posterior covariance of the process at arbitrary inputs can be expressed as a combination of a finite set of parameters which depend on the training data only. For arbitrary likelihoods we can show that

Lemma 1 (Parametrisation). *The result of the Bayesian update eq. (1) using a GP prior with mean function $\langle f_{\mathbf{x}} \rangle_0$ and kernel $K_0(\mathbf{x}, \mathbf{x}')$ and data $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n) | n = 1, \dots, N\}$ is a process with mean and kernel functions given by*

$$\begin{aligned} \langle f_{\mathbf{x}} \rangle_{\text{post}} &= \langle f_{\mathbf{x}} \rangle_0 + \sum_{i=1}^N K_0(\mathbf{x}, \mathbf{x}_i) q(i) \\ K_{\text{post}}(\mathbf{x}, \mathbf{x}') &= K_0(\mathbf{x}, \mathbf{x}') + \sum_{i,j=1}^N K_0(\mathbf{x}, \mathbf{x}_i) R(ij) K_0(\mathbf{x}_j, \mathbf{x}'). \end{aligned} \quad (2)$$

The parameters $q(i)$ and $R(ij)$ are given by

$$\begin{aligned} q(i) &= \frac{1}{Z} \int d\mathbf{f} p_0(\mathbf{f}) \frac{\partial P(\mathcal{D}|\mathbf{f})}{\partial f(\mathbf{x}_i)} \quad \text{and} \\ R(ij) &= \frac{1}{Z} \int d\mathbf{f} p_0(\mathbf{f}) \frac{\partial^2 P(\mathcal{D}|\mathbf{f})}{\partial f(\mathbf{x}_i) \partial f(\mathbf{x}_j)} - q(i)q(j) \end{aligned} \quad (3)$$

where $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^T$ and $Z = \int d\mathbf{f} p_0(\mathbf{f}) P(\mathcal{D}|\mathbf{f})$ is a normalising constant.

The parameters $q(i)$ and $R(ij)$ have to be computed only once during the training of the model, and are fixed when we make predictions. The parametric form of the posterior mean (assuming a zero mean for the prior) resembles the representations for the predictors in other kernel approaches (such as Support Vector machines) that are obtained by minimising certain cost functions. While the latter representations are derived from the celebrated representer theorem of Kimeldorf and Wahba (1971) our result (eq.2) does to our best knowledge not follow from this but is derived from simple properties of Gaussian distributions. To keep focused on the main flow, we defer the proof to Appendix B.

Making an immediate use of this representation is usually not possible because the posterior process is in general not Gaussian and the integrals cannot be computed exactly. Hence, we need approximations in order to keep the inference tractable (Csató et al. 2000). One popular method is to approximate the posterior by a Gaussian process (Williams and Barber 1998; Seeger 2000). This may be formulated within a variational approach, where a certain dissimilarity measure between the true and the approximate distribution is minimised. The most popular choice is the Kullback-Leibler divergence between distributions defined as

$$KL(p|q) = \int d\theta p(\theta) \ln \frac{p(\theta)}{q(\theta)} \quad (4)$$

where θ denotes the set of arguments of the densities. If \hat{p} denotes the approximating Gaussian distribution, one usually tries to minimise $KL(\hat{p}||p_{\text{post}})$, with respect to \hat{p} which in contrast to $KL(p_{\text{post}}||\hat{p})$ requires only the computation of expectations over tractable distributions.

with coefficients $\alpha_t(i)$ and $C_t(ij)$ not depending on \mathbf{x} and \mathbf{x}' (for details see Appendix C). For simplicity the values $\alpha_t(i)$ are grouped into the vector $\boldsymbol{\alpha}_t = [\alpha_t(1), \dots, \alpha_t(t)]^T$, $\mathbf{C}_t = \{C_t(ij)\}_{i,j=1,t}$ and we also used vectorial (typeset in bold) notations for $\mathbf{k}_x = [K_0(\mathbf{x}_1, \mathbf{x}), \dots, K_0(\mathbf{x}_t, \mathbf{x})]^T$.

The recursion for the GP parameters in eq. (8) are found from the recursion eq (6) and the parametrisation lemma:

$$\begin{aligned}\boldsymbol{\alpha}_{t+1} &= \mathbf{T}_{t+1}(\boldsymbol{\alpha}_t) + \mathbf{q}^{(t+1)} \mathbf{s}_{t+1} \\ \mathbf{C}_{t+1} &= \mathbf{U}_{t+1}(\mathbf{C}_t) + \mathbf{r}^{(t+1)} \mathbf{s}_{t+1} \mathbf{s}_{t+1}^T \\ \mathbf{s}_{t+1} &= \mathbf{T}_{t+1}(\mathbf{C}_t \mathbf{k}_{t+1}) + \mathbf{e}_{t+1}\end{aligned}\tag{9}$$

where $\mathbf{k}_{t+1} = \mathbf{k}_{\mathbf{x}_{t+1}}$ and \mathbf{e}_{t+1} the $t+1$ -th unit vector and \mathbf{s}_{t+1} is introduced for clarity. We also introduced the operators \mathbf{T}_{t+1} and \mathbf{U}_{t+1} , they extend a t -dimensional vector and matrix to a $t+1$ -dimensional one by appending zeros at the end of the vector and to the last row and column of the matrix respectively.

Since \mathbf{e}_{t+1} is the $t+1$ -th unit vector, we see that the dimension of the vector $\boldsymbol{\alpha}$ and the size of matrix \mathbf{C} increases with each likelihood point added.

Equations (6) and (7) show some resemblance to the well known extended Kalman filter. This is to be expected, because the latter approach can also be understood as a sequential propagation of an approximate Gaussian distribution. However, the main difference between the two methods is in the way the likelihood model is incorporated. While the extended Kalman filter (see Bottou (1998) for a general framework) is based on a linearisation of the likelihood, our approach uses a more robust *smoothing* of the likelihood instead.

The drawback of using (9) in practice is the quadratic increase of the number of parameters with the number of training examples. This is a feature common to most other methods of inference with Gaussian processes. A modification of the learning rule that controls the number of parameters is the main contribution of this paper and is detailed in the following.

3 Sparseness in Gaussian Processes

Sparseness can be introduced within the GP framework by using suitable approximations on the representation eq. (8). Our goal is to perform an update *without increasing* the number of parameters $\boldsymbol{\alpha}$ and \mathbf{C} when, according to a certain criteria, the error due to the approximation is not too large. This could be achieved exactly, if the new input \mathbf{x}_{t+1} would be such that the relation

$$K_0(\mathbf{x}, \mathbf{x}_{t+1}) = \sum_{i=1}^t \hat{\mathbf{e}}_{t+1}(i) K_0(\mathbf{x}, \mathbf{x}_i)\tag{10}$$

holds for all \mathbf{x} . In such a case we would have a representation for the updated process in the form eq. (8) using only the first t inputs, but with “renormalised” parameters $\hat{\boldsymbol{\alpha}}$ and $\hat{\mathbf{C}}$. A glance at eq. (9) shows that the only change would be the replacement of the vector \mathbf{s}_{t+1} by

$$\hat{\mathbf{s}}_{t+1} = \mathbf{C}_t \mathbf{k}_{t+1} + \hat{\mathbf{e}}_{t+1}.\tag{11}$$

Note, that $\hat{\mathbf{e}}_{t+1}$ is a vector of dimensionality t ! Unfortunately, for most kernels and inputs \mathbf{x}_{t+1} (10) can not be fulfilled for *all* \mathbf{x} . Nevertheless, as an approximation, we could try an update of the form (11) where $\hat{\mathbf{e}}_{t+1}$ is determined by minimising the error measure

$$\left\| K_0(\cdot, \mathbf{x}_{t+1}) - \sum_{i=1}^t \hat{\mathbf{e}}_{t+1}(i) K_0(\cdot, \mathbf{x}_i) \right\|^2,\tag{12}$$

where $\|\cdot\|$ is a suitably defined norm in a space of functions of inputs \mathbf{x} (related optimisation criteria in a function space are presented by Vijayakumar and Ogawa (1999)). Eq. (12) becomes especially simple, when the norm is based on the inner product of the *reproducing kernel Hilbert space* RKHS generated by the kernel K_0 . In this case, for any two functions g and h that are represented as $g(\mathbf{x}) = \sum_i c_i K_0(\mathbf{x}, \mathbf{u}_i)$ and $h(\mathbf{x}) = \sum_i d_i K_0(\mathbf{x}, \mathbf{v}_i)$, for some arbitrary set of \mathbf{u}_i 's and \mathbf{v}_i 's, the RKHS inner product is defined as (Wahba 1990):

$$(g(\cdot), h(\cdot))_{\text{RKHS}} = \sum_{ij} c_i d_j K_0(\mathbf{u}_i, \mathbf{v}_j)\tag{13}$$

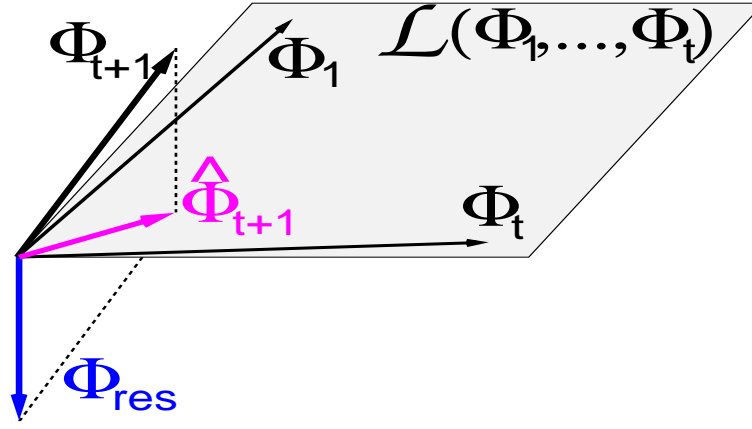


Figure 2: Visualisation of the projection step. The new feature vector ϕ_{t+1} is projected to the subspace spanned by $\{\phi_1, \dots, \phi_t\}$ resulting in the projection $\hat{\phi}_{t+1}$ and the orthogonal residual (the “left out” quantity) ϕ_{res} . It is important that ϕ_{res} has $t+1$ components, i.e. it needs the extended basis including ϕ_{t+1} .

with norm

$$\|g\|_{\text{RKHS}}^2 = (g(\cdot), g(\cdot))_{\text{RKHS}} = \sum_{ij} c_i c_j K_0(z_i, z_j). \quad (14)$$

Hence, in this case eq. (12) is

$$K_0(x_{t+1}, x_{t+1}) + \sum_{i,j=1}^t \hat{e}_{t+1}(i) \hat{e}_{t+1}(j) K_0(x_i, x_j) - 2 \sum_{i=1}^t \hat{e}_{t+1}(i) K_0(x_{t+1}, x_i) \quad (15)$$

and simple minimisation of eq. (15) yields (Smola and Schölkopf 2000)

$$\hat{e}_{t+1} = \mathbf{K}_t^{(-1)} \mathbf{k}_{t+1} \quad (16)$$

where $\mathbf{K}_t = \{K_0(x_i, x_j)\}_{i,j=1,t}$ is the Gram matrix. The expression

$$\hat{K}_0(x, x_{t+1}) = \sum_{i=1}^t \hat{e}_{t+1}(i) K_0(x, x_i) \quad (17)$$

is simply the orthogonal projection (in the sense of the inner product eq. 13) of the function $K_0(x, x_{t+1})$ on the linear span of the functions $K_0(x, x_i)$. The approximate update using (11) will be performed only, when a certain measure for the approximation error (to be discussed later) is not exceeded. The set of inputs, for which the *exact* update is performed, and the number of parameters is increased, will be called “*basis vector set*” or \mathcal{BV} set, an element will be \mathcal{BV} . Proceeding sequentially, some of the inputs are left out and others are included in \mathcal{BV} set. However, due to the projection (17) the inputs left out from \mathcal{BV} set will still contribute to the final GP configuration – the one used for prediction and to measure the posterior uncertainties. But the latter inputs *will not be stored* and do not lead to an increase of the size of the parameter set.

This procedure leads to the new representation for the posterior GP *only* in terms of the \mathcal{BV} set and the corresponding parameters α and \mathbf{C} :

$$\begin{aligned} \langle f_x \rangle &= \sum_{i \in \mathcal{BV}} K_0(x, x_i) \alpha(i) \\ K(x, x') &= K_0(x, x') + \sum_{i,j \in \mathcal{BV}} K_0(x, x_i) C(ij) K_0(x_j, x'). \end{aligned} \quad (18)$$

An alternative derivation of these results can be obtained from the representation of the Mercer kernel (Wahba 1990; Vapnik 1995)

$$K_0(x, x') = \phi(x)^T \phi(x'), \quad (19)$$

in terms of the possibly infinite dimensional “feature vector” $\phi(x)$ (Wahba 1990). Minimising (15) and using (11) for an update is equivalent to replacing the feature vector ϕ_{t+1} corresponding to the new input by its orthogonal

projection

$$\hat{\phi}_{t+1} = \sum_i \hat{e}_{t+1}(i) \phi_i \quad (20)$$

onto the space of the old feature vectors (as in Fig. 2). Note however that this derivation may be somewhat misleading, by suggesting that the mapping induced by the feature vectors plays a special role for our approximation. This would be confusing because the representation eq. (19) is not unique. Our first derivation based on the RKHS norm shows however that our approximation uses only geometrical properties that are induced by the kernel K_0 .

3.1 Projection-Induced Error

We need a rule to decide if the current input will be included in the \mathcal{BV} set or not. We base the decision on a measure of change on the sample averaged posterior mean of the GP due to the sparse approximation.

Assuming a learning scenario where only the basis vectors are memorised, we measure the change of the posterior mean due to the approximation by

$$\Delta \langle f_x \rangle_{t+1} = \langle f_x \rangle_{t+1} - \langle f_x \rangle_{\widehat{t+1}}$$

where $\langle f_x \rangle_{\widehat{t+1}}$ is the posterior mean with respect to the approximated process. Summing up the absolute values of the changes for the elements in the \mathcal{BV} set and the new input leads to

$$\begin{aligned} \varepsilon_{t+1} &= \sum_{i=1}^{t+1} |\Delta \langle f_i \rangle_{t+1}| = |q^{t+1}| \sum_{i=1}^{t+1} \left| K_0(x_i, x_{t+1}) - \widehat{K}_0(x_i, x_{t+1}) \right| \\ &= |q^{(t+1)}| \left\| K_0(\cdot, x_{t+1}) - \widehat{K}_0(\cdot, x_{t+1}) \right\|_{\text{RKHS}}^2 \end{aligned} \quad (21)$$

where the second line follows from the orthogonal projection together with the definition of the inner product in the RKHS (eq. 14)

It is an important observation that, also due to orthogonal projection, the error is concentrated only on the last data point since $\widehat{K}_0(\cdot, x_{t+1}) = K_0(\cdot, x_{t+1})$ at the old data points x_i , $i = 1, \dots, t$. Rewriting eq. (21) using the coefficients for $\widehat{K}_0(\cdot, x_{t+1})$ from eq. (16), the error is

$$\varepsilon_{t+1} = |q^{(t+1)}| \left(k_{t+1}^* - \mathbf{k}_{t+1}^T \mathbf{K}_t^{-1} \mathbf{k}_{t+1} \right) = |q^{(t+1)}| \gamma_{t+1} \quad (22)$$

where $k_{t+1}^* = K_0(x_{t+1}, x_{t+1})$ and $q^{(t+1)}$ is given from eq. (7). The error measure ε_{t+1} is a product of two terms. If the new input would be included into \mathcal{BV} the corresponding coefficient α_{t+1} in the posterior mean would be equal to $q^{(t+1)}$, which is the “likelihood-dependent” part. The second term

$$\gamma_{t+1} = k_{t+1}^* - \mathbf{k}_{t+1}^T \mathbf{K}_t^{-1} \mathbf{k}_{t+1} \quad (23)$$

gives the *geometrical part*, which is the squared norm of the “residual vector” from the projection in the RKHS (shown in Fig. 2), or equivalently the “novelty” of the current input. If we use the RBF kernels, then the error eq. (22) is similar to the one used in deciding if new centres have to be included in resource allocating network (McLachlan and Lowe 1996; Platt 1991).

To compute the geometrical component of the error ε_{t+1} , a matrix inversion is needed at each step. The costly matrix inversion can be avoided by keeping track of the inverse Gram matrix $\mathbf{Q}_t = \mathbf{K}_t^{-1}$. The updates for the matrix can also be expressed with the variables γ_{t+1} and $\hat{\mathbf{e}}_{t+1}$ (for details see D), and these updates will be important when deleting a \mathcal{BV} :

$$\mathbf{Q}_{t+1} = \mathbf{U}_{t+1}(\mathbf{Q}_t) + \gamma_{t+1}^{-1} (\mathbf{T}_{t+1}(\hat{\mathbf{e}}_{t+1}) - \mathbf{e}_{t+1})(\mathbf{T}_{t+1}(\hat{\mathbf{e}}_{t+1}) - \mathbf{e}_{t+1})^T. \quad (24)$$

where \mathbf{U}_{t+1} and \mathbf{T}_{t+1} are the extension operators for a matrix and a vector respectively (introduced in eq. (9)).

3.2 Deleting a Basis Vector

Our algorithm may run into problems, when there is no possibility to include new inputs into \mathcal{BV} without deleting one of the old basis vectors because we are at the limit of our resources. This gives the motivation to implement

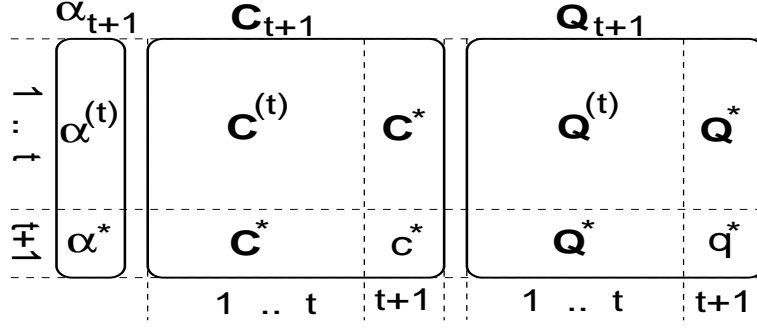


Figure 3: Grouping of the GP parameters for the update equation (25).

pruning: whenever a new example is found important, one should get rid of the basis vector (\mathcal{BV}) with the smallest error and replace it by the new input vector. First we will discuss the elimination of a \mathcal{BV} and then the criterion based on which we choose the \mathcal{BV} to be removed.

To remove a basis vector from the \mathcal{BV} set we first assume that the respective \mathcal{BV} *has just been added* – thus the previous update step was done with \mathbf{e}_{t+1} ; the $t+1$ -th unit vector. With this assumption we identify the elements $\mathbf{q}^{(t+1)}$, $\mathbf{r}^{(t+1)}$ and \mathbf{s}_{t+1} from eq. (9), compute $\hat{\mathbf{e}}_{t+1}$ (this computation is also replaced by an identification from eq. (24)) and use eq. (20) for an update without including the new point in the \mathcal{BV} set.

If we assume $t+1$ basis vectors, $\boldsymbol{\alpha}_{t+1}$ has $t+1$ elements, and the matrices \mathbf{C}_{t+1} and \mathbf{Q}_{t+1} are $(t+1) \times (t+1)$. Further assuming that we want to delete the last added element, the decomposition is as illustrated in Fig. 3. Computing the “previous” model parameters and then using the non-increasing update leads to the deletion equations (see Appendix E for details):

$$\begin{aligned}\hat{\boldsymbol{\alpha}} &= \boldsymbol{\alpha}^{(t)} - \alpha^* \frac{\mathbf{Q}^*}{q^*} \\ \hat{\mathbf{C}} &= \mathbf{C}^{(t)} + \mathbf{c}^* \frac{\mathbf{Q}^* \mathbf{Q}^{*T}}{q^{*2}} - \frac{1}{q^*} [\mathbf{Q}^* \mathbf{C}^{*T} + \mathbf{C}^* \mathbf{Q}^{*T}] \\ \hat{\mathbf{Q}} &= \mathbf{Q}^{(t)} - \frac{\mathbf{Q}^* \mathbf{Q}^{*T}}{q^*}\end{aligned}\tag{25}$$

where $\hat{\boldsymbol{\alpha}}$, $\hat{\mathbf{C}}$ and $\hat{\mathbf{Q}}$ are the parameters after the deletion of the last basis vector and $\mathbf{C}^{(t)}$, $\mathbf{Q}^{(t)}$, $\boldsymbol{\alpha}^{(t)}$, \mathbf{Q}^* , \mathbf{C}^* , \mathbf{q}^* , and \mathbf{c}^* are taken from GP parameters before deletion. A graphical illustration of each element is provided in Fig. 3.

Of particular interest is the identification of the parameters $\mathbf{q}^{(t+1)}$ and γ_{t+1} since their product gives the score of the basis vector that is being deleted. This leads to the score

$$\varepsilon_{t+1} = \frac{\alpha^*}{q^*} = \frac{\alpha_{t+1}(t+1)}{Q_{t+1}(t+1, t+1)}\tag{26}$$

Thus we have the score for the last input point. Neglecting dependence of the GP posterior on the ordering of the data, (26) gives us a score measure for each element i in the \mathcal{BV} set:

$$\varepsilon_i = \frac{|\alpha_{t+1}(i)|}{Q_{t+1}(i, i)}.\tag{27}$$

by rearranging the order in \mathcal{BV} with element i at the last position. To summarise, if a deletion is needed, then the basis vector with minimal score (from eq. (27)) will be deleted. The scores are computationally cheap (linear).

3.3 The Sparse GP Algorithm

The following numerical experiments are based on the version of the algorithm that assumes a given maximal size for the \mathcal{BV} set.

We start by initialising the \mathcal{BV} set with an empty set, the maximal number of elements in the \mathcal{BV} set with \mathbf{d} , the prior kernel \mathbf{K}_0 , and a tolerance parameter ϵ_{tol} . This latter will be used to prevent the Gram matrix from being singular and is used in step 2. The GP parameters $\boldsymbol{\alpha}$, \mathbf{C} , and the inverse Gram matrix \mathbf{Q} are set to empty values.

For each data element $(\mathbf{y}_{t+1}, \mathbf{x}_{t+1})$ we will iterate the following steps:

1. Compute $\mathbf{q}^{(t+1)}$, $\mathbf{r}^{(t+1)}$, \mathbf{k}_{t+1}^* , \mathbf{k}_{t+1} , $\hat{\mathbf{e}}_{t+1}$, and γ_{t+1} .
2. If $\gamma_{t+1} < \epsilon_{\text{tol}}$ then perform a reduced update with $\hat{\mathbf{s}}_{t+1}$ from eq. (11) without extending the size of the parameters $\boldsymbol{\alpha}$ and \mathbf{C} . Advance to the next data.
3. (else) Perform the update eq. (9) using the unit vector \mathbf{e}_{t+1} . Add the current input to the \mathcal{BV} set and compute the inverse of the extended Gram matrix using eq. (24).
4. If the size of the \mathcal{BV} set is larger than \mathbf{d} , then compute the scores ϵ_i for all \mathcal{BV} s from eq. (27), find the basis vector with the minimum score and delete it using eqs. (25).

The computational time for a single step is quadratic in \mathbf{d} , the maximal number of \mathcal{BV} s allowed. Having an iteration over all data, the computational time is $\mathcal{O}(\mathbf{N}\mathbf{d}^2)$. This is a significant improvement from the $\mathcal{O}(\mathbf{N}^3)$ scaling of the non-sparse GP algorithm. Since we propose a “subspace” algorithm, we have the same computing time as in (Wahba 1990) and the Nyström approximation for kernel matrices in (Williams and Seeger 2001).

An important feature of the sparse GP algorithm is that we provide an approximation to the *posterior kernel* of the process providing predictive variance for a new inputs, as shown by the error bars in Fig. 4. An other aspect of the algorithm is that the basis vectors are selected during runtime from the data. A final remark is that the iterative computation of the inverse Gram matrix \mathbf{Q} allows to stop from being ill-defined: γ_{t+1} is zero if the new element \mathbf{x}_{t+1} would make the Gram matrix singular. The comparison with a preset tolerance value prevents this.

4 Experimental Results

In all experiments we used spherical RBF kernels

$$\mathbf{K}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\mathbf{d}\sigma_K^2}\right) \quad (28)$$

where σ_K is the width of the kernel and \mathbf{d} is the input dimension.

4.1 Regression

In the regression model, we assume a multidimensional input $\mathbf{x} \in \mathbb{R}^m$ and an output \mathbf{y} with the likelihood

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left\{-\frac{\|\mathbf{y} - \mathbf{f}_\mathbf{x}\|^2}{2\sigma_0^2}\right\}. \quad (29)$$

Since the likelihood is Gaussian, the use of a Gaussian posterior in the on-line algorithm is *exact*. Hence, only the sparsity will introduce an approximation into our procedure. For a given number of examples, the parametrisation (8) of the posterior in terms of $\boldsymbol{\alpha}$ and \mathbf{C} leads to a predictive distribution of \mathbf{y} for an input \mathbf{x}

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \mathbf{C}) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left\{-\frac{\|\mathbf{y} - \boldsymbol{\alpha}^\top \mathbf{k}_\mathbf{x}\|^2}{2\sigma_x^2}\right\} \quad (30)$$

with $\sigma_x^2 = \sigma_0^2 + \mathbf{k}_\mathbf{x}^\top \mathbf{C}_t \mathbf{k}_\mathbf{x} + k_\mathbf{x}^*$. The online update rules eq. (9) for $\boldsymbol{\alpha}$ and \mathbf{C} in terms of the parameters $\mathbf{q}^{(t+1)}$ and $\mathbf{r}^{(t+1)}$ are:

$$\mathbf{q}^{(t+1)} = (\mathbf{y} - \boldsymbol{\alpha}_t^\top \mathbf{k}_\mathbf{x})/\sigma_x^2 \quad \mathbf{r}^{(t+1)} = -\frac{1}{\sigma_x^2}. \quad (31)$$

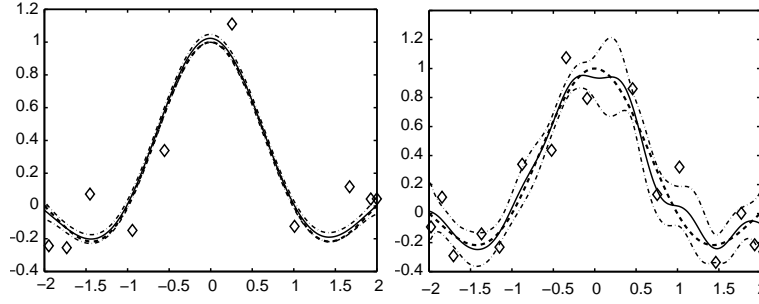


Figure 4:

Results of the GP regression with 1000 noisy training data points with noise $\sigma_0^2 = 0.02$. The figures show the results for RBF kernels with two different widths. In the left figure the good fit of the GP mean function (continuous lines) to the true function (dashed lines) is also consistent with the tight *Bayesian error-bars* (dash-dotted lines) around the means. In the right figure, the error bars are broader, reflecting the larger uncertainty. The \mathcal{BV} set is marked with rhombs and we kept 10 and 15 basis vectors. We used $\sigma_K^2 = 1$ for the left, and $\sigma_K^2 = 0.1$ for the right subfigure respectively.

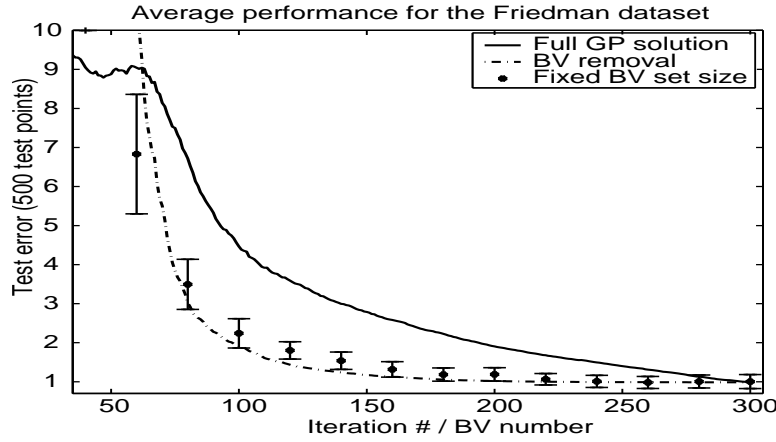


Figure 5:

Results for the Friedman data using the full GP regression (continuous line), the proposed sparse GP algorithm with a fixed \mathcal{BV} size (dots with error bars). The dash-dotted line is obtained by sequentially reducing the size of the \mathcal{BV} set. The lines show the average performance over 50 runs. The “Full GP solution” uses only the specified number of data whereas the other two curves are obtained by iterating over the full dataset ($\sigma_K^2 = 1$ was used with 300 training and 500 test data).

To illustrate the performance, we begin with the toy example $y = \sin(\pi x)/(\pi x) + \eta$ where η is a zero-mean Gaussian noise with variance $\sigma_0^2 = 0.02$. The results for the posterior means and the Bayesian error bars together with the basis vectors are shown in Fig. 4. The large error bars obtained for the “misspecified” kernel (with small width $\sigma_K^2 = 0.1$) demonstrate the advantage of propagating not only the means but also the uncertainties in the algorithm.

The second dataset is the Friedman dataset #1 (Friedman 1991), an artificial dataset frequently used to assess the performance of regression algorithms. For this example, we demonstrate the effect of the approximation introduced by the sparseness. The upper solid line in (5) shows the development of the test error with increasing numbers of examples without sparseness, i.e. when *all* data are included sequentially. The dots are the test errors obtained by running the sparse algorithm using different sizes of the \mathcal{BV} set. We see that almost two thirds of the original training set can be excluded from the \mathcal{BV} set without a significant loss of predictive performance. Finally, we have tested the effect of the greediness of the algorithm by adding or removing examples in different ways. The dependence on the data of the sparse GP is shown with the the error bars around the dots, and the dependence of result on the different orders is well within these error-bars. The dash-dotted line is obtained by first running the on-line algorithm *without sparseness* on the full data set and then building sets \mathcal{BV} of *decreasing* sizes by removing the least significant examples one after the other. Remarkably, the performance is rather stable against these variations in the plateau region of (almost) constant test error.

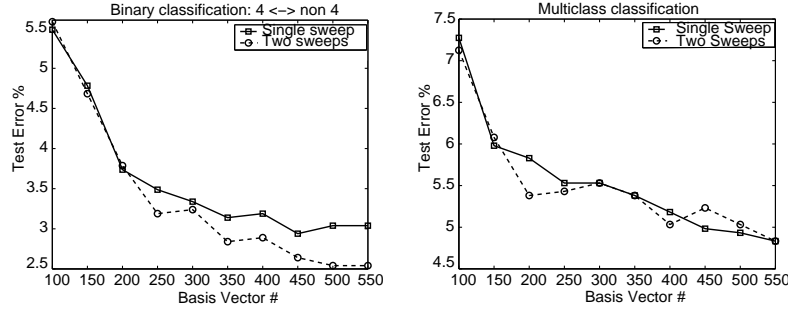


Figure 6: Results for the binary (left) and multi-class (b) classification. The multi-class case is a combination of the 10 individual classifiers: the example \mathbf{x} is assigned to the class with highest $P(C_i|\mathbf{x})$. We compare different sizes of the \mathcal{BV} set and the effect of reusing data a second (circles) time.

4.2 Classification

For classification we use the probit model (Neal 1997) where a binary value $y \in \{-1, 1\}$ is assigned to an input $\mathbf{x} \in \mathbb{R}^m$ with the data likelihood

$$P(y|\mathbf{f}_\mathbf{x}) = \text{Erf}\left(\frac{y \mathbf{f}_\mathbf{x}}{\sigma_0}\right), \quad (32)$$

$\text{Erf}(\mathbf{x})$ is the cumulative Gaussian distribution¹, with σ_0 the noise variance. The predictive distribution for a new example \mathbf{x} is:

$$p(y|\mathbf{x}, \boldsymbol{\alpha}, \mathbf{C}) = \langle P(y|\mathbf{f}_\mathbf{x}) \rangle_t = \text{Erf}\left(\frac{y \langle \mathbf{f}_\mathbf{x} \rangle}{\sigma_\mathbf{x}}\right) \quad (33)$$

where $\langle \mathbf{f}_\mathbf{x} \rangle$ the mean of the GP at \mathbf{x} given by eq. (8) and $\sigma_\mathbf{x}^2 = \sigma_0^2 + \mathbf{k}_\mathbf{x}^* + \mathbf{k}_\mathbf{x}^\top \mathbf{C} \mathbf{k}_\mathbf{x}$. Based on eq. (9), for a given input-output pair (\mathbf{x}, y) the update coefficients $q^{(t+1)}$ and $r^{(t+1)}$ are computed (for details see (Csat  et al. 2000)):

$$q^{(t+1)} = \frac{y}{\sigma_\mathbf{x}} \frac{\text{Erf}'}{\text{Erf}} \quad r^{(t+1)} = \frac{1}{\sigma_\mathbf{x}^2} \left\{ \frac{\text{Erf}''}{\text{Erf}} - \left(\frac{\text{Erf}'}{\text{Erf}} \right)^2 \right\} \quad (34)$$

with $\text{Erf}(z)$ evaluated at $z = \frac{y \boldsymbol{\alpha}_t^\top \mathbf{k}_\mathbf{x}}{\sigma_\mathbf{x}}$ and Erf' and Erf'' are the first and second derivatives at z .

We have tested the sparse GP algorithm on the USPS dataset² of gray-scale handwritten digit images (of size 16×16) with 7291 training patterns and 2007 test patterns. In the first experiment we studied the problem of classifying the digit 4 against all other digits. Fig. 6.a plots the test errors of the algorithm for different \mathcal{BV} set sizes and fixed values of hyperparameter $\sigma_k^2 = 1$.

The USPS dataset has been used previously to test the performance of other kernel-based classification algorithms that are based on a sparse representations. We mention the kernel PCA method of (Sch lkopf et al. 1999) or the Nystr m method of (Williams and Seeger 2001). They obtained slightly better results than our on-line algorithm. When the basis of the Nystr m approach is reduced to 512 the mean error is $\approx 1.7\%$ (Williams and Seeger 2001) and the PCA reduced-set method of (Sch lkopf et al. 1999) leads to an error rate of $\approx 5\%$. This may be due to the fact that the sequential replacement of the posterior by a Gaussian is an approximation for the classification problem. Hence, some of the information contained in an example is lost even when the \mathcal{BV} set would contain all data. As shown in in Fig 6 we observe a slight improvement when the algorithm sweeps several times through the data. However, it should be noted that the use of the algorithm (in its present form) on data that it has already seen is a mere heuristic and can no longer be justified from a probabilistic point of view. A change of the update rule based on a recycling of examples will be investigated elsewhere.

We have also tested our method on the more realistic problem of classifying all ten digits simultaneously. Our ability to compute Bayesian predictive probabilities is absolutely essential in this case. We have trained 10 classifiers on the ten binary classification problems of separating a single digit from the rest. A new input was assigned to the class with the highest predictive probability given by eq. (33). Fig. 6 summarises the results for the multi-class case for different \mathcal{BV} set sizes and Gaussian kernels (with the external noise variance $\sigma_0^2 = 0$). In this case, the recycling of examples

¹ $\text{Erf}(\mathbf{x}) = \int_{-\infty}^{\mathbf{x}} dt \exp(-t^2/2)/\sqrt{2\pi}$

²Available from <http://www.kernel-machines.org/data/>

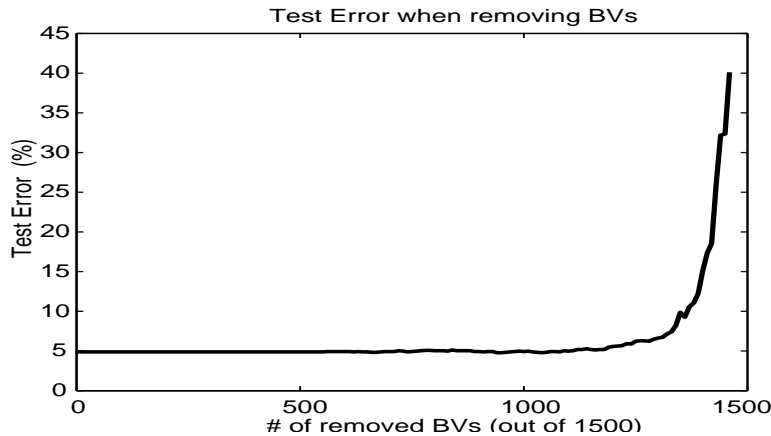


Figure 7: The performance of the combined classifier trained with an initial \mathcal{BV} size of 1500 and a sequential removal of basis vectors.

was of less significance. The gap between our on-line result and the batch performance reported in (Schölkopf et al. 1999) is also smaller, this might be due to the Bayesian nature of the GPs that avoids the over-fitting.

To reduce the computational cost we used the same set for all individual classifiers (only a single inverse of the Gram matrix was needed and also the storage cost is smaller). This made the implementation of deleting a basis vector for the multi-class case less straightforward: for each input and each basis vector there are 10 individual scores. We implemented a “minimax” deletion rule: whenever a deletion was needed, the basis vector having the smallest maximum value among the 10 classifier problems was deleted, i.e. the index l of the deleted input was

$$l = \arg \min_{i \in \mathcal{BV}} \max_{c \in \{0, \dots, 9\}} \varepsilon_i^c \quad (35)$$

Fig. 7 shows the evolution of the test error when the sparse GP algorithm was initially trained with 1500 \mathcal{BV} s and (without any retraining) the “least scoring” basis vectors are deleted. Like in the regression case (Fig. 5) we observe a long plateau of almost constant test error when up to 70% of the \mathcal{BV} s are removed.

5 Conclusions and further investigations

We have presented a greedy algorithm which allows to compute a sparse Gaussian approximation for the posterior of GP models with general (factorising) likelihoods which is based on a single sweep through the data set. So far, we have applied the method to regression and classification tasks and obtained a performance close to batch methods.

The strength of the method lies in the fact that arbitrary, even non-continuous likelihoods which maybe zero in certain regions, can be treated by our method. Such likelihoods may cause problems for other Gaussian approximations based on local linearisations (advanced Kalman filters) or on the averaging of the log-likelihood (variational Gaussian approximation). Our method merely requires the explicit computation of a Gaussian smoothed likelihood and is thus well suited for cases, where (local) likelihood functions can be modelled empirically as mixtures of Gaussians. If such expressions are available, the necessary one-dimensional integrals can be done analytically and the on-line updates require just matrix multiplications and function evaluations. A model of this structure for which we already have obtained promising preliminary results is the one used to predict wind-fields from ambiguous satellite measurements based on a GP prior for the wind-fields and a likelihood model for the measurement process.

However, a further development of the method requires the solution of various theoretical problems at which we are presently working. An important problem is to assess the quality of our approximations. There are two sources of errors. One coming from the Gaussian on-line approximation and another stemming from the additional sparsity. In both cases it is easy to obtain explicit expressions for single step errors but it is not obvious how to combine these in order to estimate the cumulative deviation between the true posterior and our approximation. It may be interesting to concentrate on the regression problem first because in this case the Gaussian approximation is exact.

A different question is the (frequentist) *statistical* quality of the algorithm. Our on-line Gaussian approximation (without sparseness) was found to be asymptotically efficient (in the sense of *Fisher*) in the finite dimensional (i.e.

parametric) case (Opper 1996; Opper 1998). This result does not trivially extend to the present infinite dimensional GP case and further studies are necessary. These may be based on the idea of an effective, finite dimensionality for the set of well estimated parameters (Trecate et al. 1999). Such work should also give an estimate for the sufficient number of basis vectors and explain the existence of the long plateaus (see Figs. 7 and 5) with practically constant test errors.

Besides a deeper understanding of the present algorithm, we find it also important to improve our method in the following ways: our sparse approximation was found to preserve the posterior means on previous data-points when projecting on a representation that leaves out the current example. A further improvement might be achieved if information on the posterior *variance* would also be used (e.g. by taking the KL loss rather than the RKHS norm) in optimising the projection. This may however result in more complex time consuming updates.

Our experiments show that in some cases the performance of the on-line algorithm is inferior to a batch method. We expect that our algorithm can be adapted to a recycling of data (e.g. along the lines of (Minka 2000)) such that a convergence to a sparse representation of the TAP mean field method (Opper and Winther 1999) is achieved.

A further drawback that will be addressed in future work is the lack of an (on-line) adaptation of the kernel hyperparameters. Rather than setting them by hand, an approximate propagation of posterior distributions for the hyperparameters would be desirable.

Finally, there may be cases of probabilistic models where the restriction to unimodal posteriors as given by the Gaussian approximation is too severe. Although we know that for Gaussian regression and classification with logistic or probit functions the posterior is unimodal, for more complicated models an on-line propagation of a mixture of GPs should be considered.

6 Acknowledgements

The authors would like to thank Bernhard Schottky for initiating the Gaussian Process parametrisation lemma. The work was supported by EPSRC grant no. GR/M81608.

References

- Berliner, L. M., L. Winkle, and N. Cressie (2000). Long-lead prediction of Pacific SST via Bayesian dynamic modelling. *Journal of Climate* 13, 3953–3968.
- Bernardo, J. M. and A. F. Smith (1994). *Bayesian Theory*. John Wiley & Sons.
- Bottou, L. (1998). Online learning and stochastic approximations. See Saad (1998), pp. 9–42.
- Cauwenberghs, G. and T. Poggio (2001). Incremental and decremental Support Vector Machine learning. In T. K. Leen, T. G. Diettrich, and V. Tresp (Eds.), *NIPS*, Volume 13. The MIT Press.
- Csató, L., E. Fokoué, M. Opper, B. Schottky, and O. Winther (2000). Efficient approaches to Gaussian process classification. In *NIPS*, Volume 12, pp. 251–257. The MIT Press.
- Evans, D. J., D. Cornford, and I. T. Nabney (2000). Structured neural network modelling for multi-valued functions for wind vector retrieval from satellite scatterometer measurements. *Neurocomputing* 30, 23–30.
- Friedman, J. H. (1991). Multivariate adaptive regression splines. *Annals of Statistics* 19, 1–141.
- Gibbs, M. and D. J. MacKay (1999). Efficient implementation of Gaussian processes. Technical report, <http://wol.ra.phy.cam.ac.uk/mackay/abstracts/gpros.html>.
- Jaakkola, T. and D. Haussler (1999). Probabilistic kernel regression. In *Online Proceedings of 7-th Int. Workshop on AI and Statistics*. <http://uncertainty99.microsoft.com/proceedings.htm>.
- Kimeldorf, G. and G. Wahba (1971). Some results on Tchebycheffian spline functions. *J. Math. Anal. Applic.* 33, 82–95.
- Lemm, J. C., J. Uhlig, and A. Weiguny (2000). A Bayesian approach to inverse quantum statistics. *Phys.Rev.Lett.* 84, 2068–20071.

- McLachlan, A. and D. Lowe (1996). Tracking of non-stationary time-series using resource allocating RBF networks. In R. Trappl (Ed.), *Cybernetics and Systems '96*, pp. 1066–1071. Proceedings of the 13th European Meeting on Cybernetics and Systems Research.
- Minka, T. P. (2000). *Expectation Propagation for Approximate Bayesian Inference*. Ph. D. thesis, Dep. of Electrical Eng. and Comp. Sci.; MIT, vismod.www.media.mit.edu/~tpminka/.
- Neal, R. M. (1997). Regression and classification using Gaussian process priors (with discussion). In J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith (Eds.), *Bayesian Statistics*, Volume 6. [ftp://ftp.cs.utoronto.ca/pub/radford/mc-gp.ps.Z](http://ftp.cs.utoronto.ca/pub/radford/mc-gp.ps.Z).
- Opper, M. (1996). Online versus offline learning from random examples: General results. *Phys. Rev. Lett.* 77(22), 4671–4674.
- Opper, M. (1998). A Bayesian approach to online learning. See Saad (1998), pp. 363–378.
- Opper, M. and O. Winther (1999). Gaussian processes and SVM: Mean field results and leave-one-out estimator. In A. Smola, P. Bartlett, B. Schölkopf, and C. Schuurmans (Eds.), *Advances in Large Margin Classifiers*, pp. 43–65. Cambridge, MA: The MIT Press.
- Platt, J. (1991). A resource-allocating network for function interpolation. *Neural Computation* 3, 213–225.
- Saad, D. (1998). *On-Line Learning in Neural Networks*. Cambridge Univ. Press.
- Schölkopf, B., S. Mika, C. J. Burges, P. Knirsch, K.-R. Müller, G. Rätsch, and A. J. Smola (1999). Input space vs. feature space in kernel-based methods. *IEEE Transactions on Neural Networks* 10(5), 1000–1017.
- Seeger, M. (2000). Bayesian model selection for Support Vector Machines, Gaussian processes and other kernel classifiers. In S. A. Solla, T. K. Leen, and K.-R. Müller (Eds.), *NIPS*, Volume 12. The MIT Press.
- Smola, A. J. and B. Schölkopf (2000). Sparse greedy matrix approximation for machine learning. In *International Conference on Machine Learning*. www.kernel-machines.org/papers/upload_4467_kfa_long.ps.gz.
- Tipping, M. (2000). The Relevance Vector Machine. In S. A. Solla, T. K. Leen, and K.-R. Müller (Eds.), *NIPS*, Volume 12. The MIT Press.
- Treccate, G. F., C. K. I. Williams, and M. Opper (1999). Finite-dimensional approximation of Gaussian processes. In M. S. Kearns, S. A. Solla, and D. A. Cohn (Eds.), *NIPS*, Volume 11. The MIT Press.
- Tresp, V. (2000). A Bayesian committee machine. *Neural Computation* 12(11), 2719–2741.
- Vapnik, V. N. (1995). *The Nature of Statistical Learning Theory*. New York, NY: Springer-Verlag.
- Vijayakumar, S. and H. Ogawa (1999). RKHS based functional analysis for exact incremental learning. *Neurocomputing* 29(1–3), 85–113.
- Wahba, G. (1990). *Splines Models for Observational Data*. Philadelphia: Series in Applied Mathematics, Vol. 59, SIAM.
- Williams, C. K. I. (1999). Prediction with Gaussian processes. In M. I. Jordan (Ed.), *Learning in Graphical Models*. The MIT Press.
- Williams, C. K. I. and D. Barber (1998). Bayesian classification with Gaussian Processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20(12), 1342–1351.
- Williams, C. K. I. and C. E. Rasmussen (1996). Gaussian processes for regression. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo (Eds.), *NIPS*, Volume 8. The MIT Press.
- Williams, C. K. I. and M. Seeger (2001). Using the Nyström method to speed up kernel machines. In T. K. Leen, T. G. Diettrich, and V. Tresp (Eds.), *NIPS*, Volume 13.

A Properties of zero-mean Gaussians

The following property of the Gaussian pdfs is often used in this paper, here we state it in a form of a theorem:

Theorem 1. Let $\mathbf{x} \in \mathbb{R}^m$ and $p(\mathbf{x})$ zero-mean Gaussian pdf with covariance $\Sigma = \{\Sigma_{ij}\}$ (i, j from 1 to m). If $g: \mathbb{R}^m \rightarrow \mathbb{R}$ is a differentiable function not growing faster than a polynomial and with partial derivatives

$$\partial_i g(\mathbf{x}) = \frac{\partial}{\partial x_i} g(\mathbf{x}) ,$$

then

$$\int_{\mathbb{R}^m} \mathbf{x} \mathbf{p}(\mathbf{x}) \mathbf{x}_i g(\mathbf{x}) = \sum_{j=1}^m \Sigma_{ij} \int_{\mathbb{R}^m} \mathbf{x} \mathbf{p}(\mathbf{x}) \partial_j g(\mathbf{x}) . \quad (36)$$

In the following we will assume definite integration over \mathbb{R}^m whenever the integral appears. Alternatively, using the vector notation, the above identity reads:

$$\int \mathbf{x} \mathbf{p}(\mathbf{x}) g(\mathbf{x}) = \mathbf{\Sigma} \int \mathbf{x} \mathbf{p}(\mathbf{x}) \nabla g(\mathbf{x}) \quad (37)$$

For a general Gaussian pdf with mean $\boldsymbol{\mu}$ the above equation transforms to:

$$\int \mathbf{x} \mathbf{p}(\mathbf{x}) g(\mathbf{x}) = \boldsymbol{\mu} \int \mathbf{x} \mathbf{p}(\mathbf{x}) g(\mathbf{x}) + \mathbf{\Sigma} \int \mathbf{x} \mathbf{p}(\mathbf{x}) \nabla g(\mathbf{x}) \quad (38)$$

Proof. The proof uses the partial integration rule:

$$\int \mathbf{x} \mathbf{p}(\mathbf{x}) \nabla g(\mathbf{x}) = - \int \mathbf{x} g(\mathbf{x}) \nabla \mathbf{p}(\mathbf{x})$$

where we have used the fast decay of the Gaussian function to dismiss one of the terms. Using the derivative of a Gaussian pdf. we have:

$$\int \mathbf{x} \mathbf{p}(\mathbf{x}) \nabla g(\mathbf{x}) = \int \mathbf{x} g(\mathbf{x}) \Sigma^{-1} \mathbf{x} \mathbf{p}(\mathbf{x})$$

Multiplying both sides with $\mathbf{\Sigma}$ leads to eq. (37), completing the proof. For the nonzero mean the deductions are also straightforward. \square

B Proof of the Parametrisation Lemma

Using Bayes' rule, the posterior process has the form

$$\hat{\mathbf{p}}(\mathbf{f}) = \frac{\mathbf{p}_0(\mathbf{f}) P(\mathcal{D}|\mathbf{f})}{\int d\mathbf{f} \mathbf{p}_0(\mathbf{f}) P(\mathcal{D}|\mathbf{f})}$$

where \mathbf{f} is a set of realisations for the random process indexed by arbitrary points from \mathbb{R}^m , the inputs for the GPs.

We compute first the mean function of the posterior process:

$$\begin{aligned} \langle f_x \rangle_{\text{post}} &= \int d\mathbf{f} \hat{\mathbf{p}}(\mathbf{f}) f_x = \frac{\int d\mathbf{f} \mathbf{p}_0(\mathbf{f}) f_x P(\mathcal{D}|\mathbf{f})}{\int d\mathbf{f} \mathbf{p}_0(\mathbf{f}) P(\mathcal{D}|\mathbf{f}) \mathbf{p}_0(\mathbf{f})} \\ &= \frac{1}{Z} \int d\mathbf{f}_x \prod_{i=1}^N d\mathbf{f}_i \mathbf{p}_0(f_x, f_1, \dots, f_N) f_x P(\mathcal{D}|f_1, \dots, f_N) \end{aligned} \quad (39)$$

where the denominator was denoted by Z and we used index notation for the realisations of the process also (thus $f(\mathbf{x}) = f_x$ and $f(\mathbf{x}_i) = f_i$). Observe that, irrespectively of the number of the random variables of the process considered, the dimension of the integral we need to consider is only $N+1$, all other random variables will integrate out (as in eq. (39)). We thus have an $N+1$ -dimensional integral in the numerator and Z is an N -dimensional integral. If we group the variables related to the data as $\mathbf{f}_D = [f_1, \dots, f_N]^T$, and apply Th. 1 (eq. 36) replacing \mathbf{x} by f_x and $g(\mathbf{x})$ by $P(\mathcal{D}|\mathbf{f}_D)$, we have

$$\begin{aligned} \langle f_x \rangle_{\text{post}} &= \frac{1}{Z} \left(\langle f_x \rangle_0 \int d\mathbf{f}_x d\mathbf{f}_D \mathbf{p}_0(f_x, \mathbf{f}_D) P(\mathcal{D}|\mathbf{f}_D) \right. \\ &\quad \left. + \sum_{i=1}^N K_0(\mathbf{x}, \mathbf{x}_i) \int d\mathbf{f}_x d\mathbf{f}_D \mathbf{p}_0(f_x, \mathbf{f}_D) \partial_i P(\mathcal{D}|\mathbf{f}_D) \right) \end{aligned} \quad (40)$$

where K_0 is the kernel function generating the covariance matrix ($\mathbf{\Sigma}$ in Theorem 1). The variable f_x in the integrals disappears since it is only contained in \mathbf{p}_0 . Substituting back Z leads to

$$\langle f_x \rangle_{\text{post}} = \langle f_x \rangle_0 + \sum_{i=1}^N K_0(\mathbf{x}, \mathbf{x}_i) q_i \quad (41)$$

where q_i is read off from eq. (40)

$$q_i = \frac{\int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) \partial_i P(\mathcal{D}|\mathbf{f}_{\mathcal{D}})}{\int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) P(\mathcal{D}|\mathbf{f}_{\mathcal{D}})} \quad (42)$$

and the coefficients q_i depend only on the data, and are independent from \mathbf{x} at which the posterior mean is evaluated.

We can simplify the expression for q_i by performing a change of variables in the numerator: $f'_i = f_i - \langle f_i \rangle_0$ where $\langle f_i \rangle_0$ is the prior mean at x_i and keeping all other variables unchanged $f'_j = f_j, j \neq i$, leading to the numerator

$$\int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}'_{\mathcal{D}}) \partial_i P(\mathcal{D}|f'_1, \dots, f'_i + \langle f_i \rangle_0, \dots, f'_N)$$

and the differentiation is with respect to f'_i . We then change the partial differentiation with respect to f'_i with the partial differentiation with respect to $\langle f_i \rangle_0$ and exchange the differentiation and integral operators (they apply to a distinct set of variables), leading to

$$\frac{\partial}{\partial \langle f_i \rangle_0} \int d\mathbf{f}'_{\mathcal{D}} p_0(\mathbf{f}'_{\mathcal{D}}) P(\mathcal{D}|f'_1, \dots, f'_i + \langle f_i \rangle_0, \dots, f'_N)$$

We then perform the inverse change of variables inside the integral and substitute back into the expression for q_i

$$q_i = \frac{\frac{\partial}{\partial \langle f_i \rangle_0} \int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) P(\mathcal{D}|\mathbf{f}_{\mathcal{D}})}{\int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) P(\mathcal{D}|\mathbf{f}_{\mathcal{D}})} = \frac{\partial}{\partial \langle f_i \rangle_0} \ln \int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) . \quad (43)$$

Writing the expression for the posterior kernel:

$$K_{\text{post}}(\mathbf{x}, \mathbf{x}') = \langle f_{\mathbf{x}} f_{\mathbf{x}'} \rangle_{\text{post}} - \langle f_{\mathbf{x}} \rangle_{\text{post}} \langle f_{\mathbf{x}'} \rangle_{\text{post}} \quad (44)$$

and applying Theorem 1 twice leads to

$$K_{\text{post}}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') + \sum_{i=1}^N \sum_{j=1}^N K_0(\mathbf{x}, \mathbf{x}_i) (D_{ij} - q_i q_j) K_0(\mathbf{x}_j, \mathbf{x}') \quad (45)$$

where D_{ij} is

$$D_{ij} = \frac{1}{Z} \int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) \frac{\partial^2}{\partial f_j \partial f_i} P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) \quad (46)$$

Identifying $R_{ij} = D_{ij} - q_i q_j$ leads to the required parametrisation in equation (3) from Lemma 1. Simplification of $R_{ij} = D_{ij} - q_i q_j$ is made by changing the arguments of the partial derivative and using the logarithm of the expectation (repeating steps (42)–(43) from q_i), leading to

$$R_{ij} = \frac{\partial^2}{\partial \langle f_i \rangle_0 \partial \langle f_j \rangle_0} \ln \int d\mathbf{f}_{\mathcal{D}} p_0(\mathbf{f}_{\mathcal{D}}) P(\mathcal{D}|\mathbf{f}_{\mathcal{D}}) \quad \square \quad (47)$$

and using a single data in the likelihood leads to the scalar coefficients $q^{(t+1)}$ and $r^{(t+1)}$ from eq. (7)

C Online Learning in GP Framework

We prove eq. (8) by induction. We will show that for every time-step we can express the mean and kernel functions with coefficients $\boldsymbol{\alpha}$ and \mathbf{C} given by the recursion (also eq. (9)):

$$\boldsymbol{\alpha}_{t+1} = T_{t+1}(\boldsymbol{\alpha}_t) + q^{(t+1)} \mathbf{s}_{t+1} \quad (48)$$

$$\mathbf{C}_{t+1} = U_{t+1}(\mathbf{C}_t) + r^{(t+1)} \mathbf{s}_{t+1} \mathbf{s}_{t+1}^T \quad (49)$$

$$\mathbf{s}_{t+1} = T_{t+1}(\mathbf{C}_t \mathbf{k}_{t+1}) + \mathbf{e}_{t+1}$$

where $\boldsymbol{\alpha}$ and \mathbf{C} depend only on the data points \mathbf{x}_i and kernel function K_0 but do not depend on the values \mathbf{x} and \mathbf{x}' (from eq. (8)) at which the mean and kernel functions are computed.

Proceeding by induction and using the induction hypothesis $\boldsymbol{\alpha}_0 = \mathbf{C}_0 = \mathbf{0}$ for time $t = 1$, we have $\alpha_1(1) = q^{(1)}$ and $C_1(1, 1) = r^{(1)}$. The mean function at time $t = 1$ has is $\langle f_{\mathbf{x}} \rangle = \alpha_1(1) K_0(\mathbf{x}_1, \mathbf{x})$ (from lemma 1 for a single data,

eq. (6)). Similarly the modified kernel is $K_1(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') + K(\mathbf{x}, \mathbf{x}_1)C_1(1, 1)K_0(\mathbf{x}_1, \mathbf{x}')$ with $\boldsymbol{\alpha}$ and \mathbf{C} independent of \mathbf{x} and \mathbf{x}' , this proving the induction hypothesis.

We assume that at time moment t we have the parameters $\boldsymbol{\alpha}_t$ and \mathbf{C}_t independent of the points \mathbf{x} and \mathbf{x}' . These parameters specify a prior GP for which we apply the online learning:

$$\begin{aligned} \langle f_{\mathbf{x}} \rangle_{t+1} &= \sum_{i=1}^t K_0(\mathbf{x}_i, \mathbf{x})\alpha_t(i) + q^{(t+1)} \sum_{i,j=1}^t K_0(\mathbf{x}, \mathbf{x}_i)C_t(i, j)K_0(\mathbf{x}_j, \mathbf{x}_{t+1}) \\ &\quad + q^{(t+1)}K_0(\mathbf{x}, \mathbf{x}_{t+1}) \\ &= \sum_{i=1}^{t+1} K_0(\mathbf{x}, \mathbf{x}_i)\alpha_{t+1}(i) \end{aligned} \quad (50)$$

and by pairwise identification we have eq. (48) or eq. (9) from the main body. The parameters $\boldsymbol{\alpha}_{t+1}$ do not depend on the particular value of \mathbf{x} . Writing down the update equation for the kernels

$$K_{t+1}(\mathbf{x}, \mathbf{x}') = K_t(\mathbf{x}, \mathbf{x}') + r^{(t+1)}K_t(\mathbf{x}, \mathbf{x}_{t+1})K_t(\mathbf{x}_{t+1}, \mathbf{x}')$$

leads to eq. (49) in a straightforward manner with $C_{t+1}(i, j)$ independent of \mathbf{x} and \mathbf{x}' , completing the induction. \square

D Iterative computation of the inverse Gram matrix

In the sparse approximation eq. (16) we need the inverse Gram matrix of the \mathcal{BV} set: $\mathbf{K}_{\mathcal{BV}} = \{K_0(\mathbf{x}_i, \mathbf{x}_j)\}$ is needed. In the following the elements of the \mathcal{BV} set are indexed from 1 to t . Using the matrix inversion formula³ the addition of a new element can be carried out sequentially. This is a well known fact, exploited also in the Kalman filter algorithm. We consider the new element at the end (last row and column) of matrix \mathbf{K}_{t+1} . Matrix \mathbf{K}_{t+1} is decomposed:

$$\mathbf{K}_{t+1} = \begin{bmatrix} \mathbf{K}_t & \mathbf{k}_{t+1} \\ \mathbf{k}_{t+1}^T & k_{t+1}^* \end{bmatrix} \quad (51)$$

Assuming \mathbf{K}_t^{-1} known and applying the matrix inversion lemma for \mathbf{K}_{t+1} :

$$\begin{aligned} \mathbf{K}_{t+1}^{-1} &= \begin{bmatrix} \mathbf{K}_t & \mathbf{k}_{t+1} \\ \mathbf{k}_{t+1}^T & k_{t+1}^* \end{bmatrix}^{-1} \\ &= \begin{bmatrix} \mathbf{K}_t^{-1} + \mathbf{K}_t^{-1}\mathbf{k}_{t+1}\mathbf{k}_{t+1}^T\mathbf{K}_t^{-1}\gamma_{t+1}^{-1} & -\mathbf{K}_t^{-1}\mathbf{k}_{t+1}\gamma_{t+1}^{-1} \\ -\mathbf{k}_{t+1}^T\mathbf{K}_t^{-1}\gamma_{t+1}^{-1} & \gamma_{t+1}^{-1} \end{bmatrix} \end{aligned} \quad (52)$$

where $\gamma_{t+1} = k_{t+1}^* - \mathbf{k}_{t+1}^T\mathbf{K}_t^{-1}\mathbf{k}_{t+1}$ is the geometric term from eq. (23). Using notations $\mathbf{K}_{t+1}\mathbf{k}_{t+1} = \hat{\mathbf{e}}_{t+1}$ from eq. (16), $\mathbf{K}_t^{-1} = \mathbf{Q}_t$, and $\mathbf{K}_{t+1}^{-1} = \mathbf{Q}_{t+1}$ we have a recursion equation:

$$\mathbf{Q}_{t+1} = \begin{bmatrix} \mathbf{Q}_t + \gamma_{t+1}^{-1}\hat{\mathbf{e}}_{t+1}\hat{\mathbf{e}}_{t+1}^T & -\gamma_{t+1}^{-1}\hat{\mathbf{e}}_{t+1} \\ -\gamma_{t+1}^{-1}\hat{\mathbf{e}}_{t+1}^T & \gamma_{t+1}^{-1} \end{bmatrix} \quad (53)$$

and in a more compact matrix notation:

$$\mathbf{Q}_{t+1} = \mathbf{Q}_t + \gamma_{t+1}^{-1}(\hat{\mathbf{e}}_{t+1} - \mathbf{e}_{t+1})(\hat{\mathbf{e}}_{t+1} - \mathbf{e}_{t+1})^T \quad (54)$$

where \mathbf{e}_{t+1} is the $t+1$ -th unit vector. With this recursion equation all matrix inversion is eliminated (this result is general for block matrices, such implementation, together with an interpretation of the parameters has been also made in (Cauwenberghs and Poggio 2001)). Using the score (26) and including in \mathcal{BV} only inputs with nonzero scores, the Gram matrix is guaranteed to be nonsingular, $\gamma_{t+1} > 0$ guarantees non-singularity of the extended Gram matrix (see Fig. 2.b).

For numerical stability we can use the Cholesky-decomposition of the inverse Gram matrix \mathbf{Q} . Using the lower-triangular matrix \mathbf{R} with the corresponding indices, and the identity $\mathbf{Q} = \mathbf{R}^T\mathbf{R}$, we have the update for the Cholesky-decomposition

$$\mathbf{R}_{t+1} = \begin{pmatrix} \mathbf{R}_t & 0 \\ -\gamma_{t+1}^{-1/2}\hat{\mathbf{e}}_{t+1}^T & \gamma_{t+1}^{-1/2} \end{pmatrix} \quad (55)$$

³A useful guide to formulae for matrix inversions and block matrix manipulation can be found at Sam Roweis' home-page: <http://www.gatsby.ucl.ac.uk/~roweis/notes.html>

that is a computationally very inexpensive operation, without additional operations provided that the quantities γ_{t+1} and \mathbf{e}_{t+1} are already computed.

E Deleting a \mathcal{BV}

Adding a basis vector is made with the equations:

$$\begin{aligned}\mathbf{s}_{t+1} &= \mathbf{T}_{t+1}(\mathbf{C}_t \mathbf{k}_{t+1}) + \mathbf{e}_{t+1} \\ \boldsymbol{\alpha}_{t+1} &= \mathbf{T}_{t+1}(\boldsymbol{\alpha}_t) + q^{(t+1)} \mathbf{s}_{t+1}\end{aligned}\tag{56}$$

$$\mathbf{C}_{t+1} = \mathbf{U}_{t+1}(\mathbf{C}_t) + r^{(t+1)} \mathbf{s}_{t+1} \mathbf{s}_{t+1}^T \tag{57}$$

$$\mathbf{Q}_{t+1} = \mathbf{U}_{t+1}(\mathbf{Q}_t) + \gamma_{t+1}^{-1} (\mathbf{T}_{t+1}(\hat{\mathbf{e}}_{t+1}) - \mathbf{e}_{t+1})(\mathbf{T}_{t+1}(\hat{\mathbf{e}}_{t+1}) - \mathbf{e}_{t+1})^T \tag{58}$$

where $\boldsymbol{\alpha}$ and \mathbf{C} are the GP parameters, \mathbf{Q} is the inverse Gram matrix, γ_{t+1} and $\hat{\mathbf{e}}_{t+1}$ the geometric terms of the new basis vector, $\mathbf{k}_{t+1} = [K_0(\mathbf{x}_1, \mathbf{x}_{t+1}), \dots, K_0(\mathbf{x}_t, \mathbf{x}_{t+1})]^T$, and \mathbf{e}_{t+1} is the $t+1$ -th unity vector.

The optimal decrease of \mathcal{BV} s needs an answer to two questions. The first question is how to delete a basis vector from the set of basis vectors with minimal loss of information. If the method is given, then we have to *find the \mathcal{BV} to remove*. The first problem is solved by “inverting” the learning equations (56–58). Assuming $\boldsymbol{\alpha}_{t+1}$, \mathbf{C}_{t+1} , and \mathbf{Q}_{t+1} known and using pairwise correspondence considering the $t+1$ -th element of $\boldsymbol{\alpha}_{t+1}$, we can identify $q^{(t+1)} = \alpha_{t+1}(t+1) \stackrel{\text{def}}{=} \alpha_{t+1}^*$ (the notations are illustrated in Fig 3). Using similar correspondences for the matrix \mathbf{C}_{t+1} the following identifications can be done:

$$r^{(t+1)} = \mathbf{C}_{t+1}(t+1, t+1) \stackrel{\text{def}}{=} c_{t+1}^* \tag{59}$$

$$\mathbf{C}_t \mathbf{k}_{t+1} = \mathbf{C}_{t+1}(1..t, t+1) \stackrel{\text{def}}{=} \frac{\mathbf{C}_{t+1}^*}{c_{t+1}^*}$$

with c_{t+1}^* and \mathbf{C}_{t+1}^* sketched in Fig. 3. Substituting back into equations (56) and (57), the old values for GP parameters are:

$$\boldsymbol{\alpha}_t = \boldsymbol{\alpha}_{t+1}^{(t)} - \alpha_{t+1}^* \frac{\mathbf{C}_{t+1}^*}{c_{t+1}^*} \tag{60}$$

$$\mathbf{C}_t = \mathbf{C}_{t+1}^{(t)} - \frac{\mathbf{C}_{t+1}^* \mathbf{C}_{t+1}^{*T}}{c_{t+1}^*} \tag{61}$$

where $\boldsymbol{\alpha}_{t+1}^{(t)} = \mathbf{T}_{t+1}^{-1}(\boldsymbol{\alpha}_{t+1})$ and \mathbf{T}_{t+1}^{-1} is the inverse operator that takes the first t elements of a $t+1$ -dimensional vector. We define $\mathbf{C}_{t+1}^{(t)} = \mathbf{U}_{t+1}^{-1}(\mathbf{C}_{t+1})$ similarly.

Proceeding similarly, using elements of matrix \mathbf{Q}_{t+1} , the correspondence with (58) is as follows:

$$\gamma_{t+1} = \frac{1}{\mathbf{Q}_{t+1}(t+1, t+1)} \stackrel{\text{def}}{=} \frac{1}{q_{t+1}^*} \tag{62}$$

$$\hat{\mathbf{e}}_{t+1} = -\frac{\mathbf{Q}_{t+1}(1..t, t+1)}{q_{t+1}^*} \stackrel{\text{def}}{=} -\frac{\mathbf{Q}_{t+1}^*}{q_{t+1}^*}$$

with the reduced set matrix $\hat{\mathbf{Q}}_{t+1}$:

$$\hat{\mathbf{Q}}_{t+1} = \mathbf{Q}_{t+1}^{(t)} - \frac{\mathbf{Q}_{t+1}^* \mathbf{Q}_{t+1}^{*T}}{q_{t+1}^*} \tag{63}$$

The matrix \mathbf{Q} does not need any further modification, however for $\boldsymbol{\alpha}$ and \mathbf{C} a sparse update is needed. Replacing γ_{t+1} , $\hat{\mathbf{e}}_{t+1}$ together with the “old” parameters $\boldsymbol{\alpha}_t$ and \mathbf{C}_t , we have the “optimally reduced” GP parameters as

$$\hat{\boldsymbol{\alpha}}_t = \boldsymbol{\alpha}_{t+1}^{(t)} + \alpha_{t+1}^* \frac{\mathbf{Q}_{t+1}^*}{q_{t+1}^*} \tag{64}$$

$$\hat{\mathbf{C}}_t = \mathbf{C}_{t+1}^{(t)} + c_{t+1}^* \frac{\mathbf{Q}_{t+1}^* \mathbf{Q}_{t+1}^{*T}}{q_{t+1}^{*2}} - \frac{1}{q_{t+1}^*} [\mathbf{Q}_{t+1}^* \mathbf{C}_{t+1}^{*T} + \mathbf{C}_{t+1}^* \mathbf{Q}_{t+1}^{*T}] \quad \square \tag{65}$$